Let

- $f_\theta$ be a language model.

- $T$ be the tokenizer associated with $f_\theta$.

- $T(w)$ denote the sequence of token IDs of a word $w$.

- $T(s)$ denote the sequence of token IDs of a sentence $s$.

- $\log \operatorname{softmax}(\mathbf{z})$ denote the vector of log-probabilities over the vocabulary obtained from logits $\mathbf{z}$.

- $\operatorname{LogSumExp}(a_1, \ldots, a_K)$ denote $\log\left(\sum_{k=1}^{K} e^{a_k}\right)$.

**Algorithm 1** , dynamic case

---

**Require:** sentence $s$, model $f_\theta$, tokenizer $T$, set of words $\mathcal{W}$

**Ensure:** mapping Result $: \mathcal{W} \to \mathbb{R}$ of average log-probabilities

1: $x \leftarrow T(s)$          $\triangleright$ $x = (x_1, \ldots, x_n)$, sentence token IDs
2: $n \leftarrow |x|$          $\triangleright$ number of tokens in the sentence
3: Result $\leftarrow \emptyset$
4: **for** each word $w \in \mathcal{W}$ **do**
5:      $\tau(w) \leftarrow T(w)$        $\triangleright$ $\tau(w) = (t_1, \ldots, t_L)$, token IDs of $w$
6:      **if** $|\tau(w)| = 0$ **then**
7:          Result$[w] \leftarrow$ NaN
8:          **continue** to next $w$
9:      $\mathcal{S} \leftarrow \emptyset$      $\triangleright$ list of joint log-probabilities for each insertion position
10:      $K \leftarrow n + 1$        $\triangleright$ number of insertion positions
11:      **for** $k = 0$ to $n$ **do**
12:          prefix $\leftarrow (x_1, \ldots, x_k)$
13:          context $\leftarrow$ prefix
14:          $\ell_{\text{joint}} \leftarrow 0$      $\triangleright$ joint log-probability for this insertion position
15:          **for** $j = 1$ to $L$ **do**
16:              **if** $|\text{context}| = 0$ **then**
17:                 eval_context $\leftarrow$ (BOS)      $\triangleright$ BOS token ID from tokenizer
18:              **else**
19:                 eval_context $\leftarrow$ context
20:              $\mathbf{z} \leftarrow f_\theta(\text{eval\_context})$        $\triangleright$ logits for next token
21:              $\ell \leftarrow \log \text{softmax}(\mathbf{z}_{\text{last}})$
22:              $\ell_{\text{token}} \leftarrow \ell[t_j]$
23:              $\ell_{\text{joint}} \leftarrow \ell_{\text{joint}} + \ell_{\text{token}}$
24:              context $\leftarrow$ context $\| t_j$        $\triangleright$ append $t_j$ to context
25:          append $\ell_{\text{joint}}$ to $\mathcal{S}$
26:      **if** $|\mathcal{S}| > 0$ **then**
27:          $\ell_{\text{sum}} \leftarrow \text{LogSumExp}(\mathcal{S})$
28:          $\ell_{\text{avg}} \leftarrow \ell_{\text{sum}} - \log K$
29:          Result$[w] \leftarrow \ell_{\text{avg}}$
30:      **else**
31:          Result$[w] \leftarrow$ NaN
32: **return** Result

---

**Algorithm 2** , static case

---

**Require:** sentence $s$, model $f_\theta$, tokenizer $T$, set of words $\mathcal{W}$
**Ensure:** mapping Result $: \mathcal{W} \to \mathbb{R}$ of average log-probabilities
1:   $x \leftarrow T(s)$                      $\triangleright$ $x = (x_1, \ldots, x_n)$, sentence token IDs
2:   $n \leftarrow |x|$
3:   $\mathbf{Z} \leftarrow f_\theta(x)$                $\triangleright$ $\mathbf{Z} \in \mathbb{R}^{n \times |\mathcal{V}|}$, logits over vocabulary $\mathcal{V}$
4:   $\mathbf{L} \leftarrow \log \mathrm{softmax}(\mathbf{Z})$ along the vocabulary dimension     $\triangleright$ $\mathbf{L} \in \mathbb{R}^{n \times |\mathcal{V}|}$ is the static log-probability table
5:   Result $\leftarrow \emptyset$
6:   **for** each word $w \in \mathcal{W}$ **do**
7:      $\tau(w) \leftarrow T(w)$               $\triangleright$ $\tau(w) = (t_1, \ldots, t_L)$, token IDs of $w$
8:      $L_w \leftarrow |\tau(w)|$
9:      **if** $L_w = 0$ **then**
10:         Result$[w] \leftarrow$ NaN
11:         **continue** to next $w$
12:      $\mathcal{S} \leftarrow \emptyset$            $\triangleright$ list of pseudo-joint log-probabilities
13:      $K \leftarrow n - L_w + 1$        $\triangleright$ number of sliding-window positions
14:      **for** $k = 1$ to $K$ **do**
15:         $\ell_{\mathrm{joint}} \leftarrow 0$
16:         **for** $j = 1$ to $L_w$ **do**
17:            $i \leftarrow k + j - 1$                $\triangleright$ sentence position
18:            $\ell_{\mathrm{token}} \leftarrow \mathbf{L}_{i, t_j}$        $\triangleright$ log-probability of $t_j$ at position $i$
19:            $\ell_{\mathrm{joint}} \leftarrow \ell_{\mathrm{joint}} + \ell_{\mathrm{token}}$
20:         append $\ell_{\mathrm{joint}}$ to $\mathcal{S}$
21:      **if** $|\mathcal{S}| > 0$ **then**
22:         $\ell_{\mathrm{sum}} \leftarrow \mathrm{LogSumExp}(\mathcal{S})$
23:         $\ell_{\mathrm{avg}} \leftarrow \ell_{\mathrm{sum}} - \log K$
24:         Result$[w] \leftarrow \ell_{\mathrm{avg}}$
25:      **else**
26:         Result$[w] \leftarrow$ NaN
27: **return** Result

---